



glushkov automata + xml

1980

- 2

Scholar All articles Recent articles Results 1 - 10 of about 24 for glushkov :

All Results

[H Hosoya](#)[B Pierce](#)[R Carrasco](#)[B Bouchou](#)[J Rico-Juan](#)

[PS] [XDuce: A Typed XML Processing Language \(Preliminary Report\)](#) - group of 9 »

H Hosoya, BC Pierce - Selected papers from the Third International Workshop WebDB, 2000 - it-c.dk

... along similar lines is the functional language XM for

XML processing, proposed ... be

close to Haskell's, except that they incorporate **Glushkov automata** in type ...

[Cited by 82](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#) - [BL Direct](#)

→ [A similarity between probabilistic tree languages: application to XML document families](#) - group of 5 »

RC Carrasco, JR Rico-Juan - Pattern Recognition, 2002 - dlsi.ua.es

... The model has been used to compute similarities between XML document sets ... states

and probabilities were extracted from the **Glushkov automata** [1] of the regular ...

[Cited by 12](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

[Handling syntactic constraints in a DTD-compliant XML editor](#) - group of 5 »

YS Kuo, J Wang, NC Shih - Proceedings of the 2003 ACM symposium on Document ..., 2003 - portal.acm.org

... To demonstrate the effectiveness of the algorithms and editing process, we build

an XML editor with forms as its user interface. 2.

GLUSHKOV AUTOMATA [1] A ...

[Cited by 4](#) - [Related Articles](#) - [Web Search](#)

→ Extending tree **automata** to model **XML** validation under element and attribute constraints - group of 5 »

B Bouchou, D Duarte, MHF Alves, D Laurent - ICEIS (1), 2003 - www-lih.univ-lehavre.fr

... fr Keywords: Regular tree **automata**, **XML** validation, **XML** views, DTD, attribute constraints, element constraints Abstract: Algorithms ... Cited by 10 - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

Automatic DTD simplification by examples

A Bia, RC Carrasco - ACH/ALLC, 2001 - nyu.edu

... the general DTD is processed to extract the structure of the markup model with which we build a **Glushkov automata** (Caron and Ziadi, 2000). The **XML** sample files ...

Cited by 5 - [Related Articles](#) - [Cached](#) - [Web Search](#)

Turning DTDs into specialized tree-**automata**-based schemata to match a collection of marked-up ... - group of 6 »

RC Carrasco, A Bia, ML Forcada, PM Perez-Anton - Rapport technique, Universidad de Alicante, 2002 - cs.technion.ac.il

... to extract the structure of the markup models and a **Glushkov automaton** [11] is built for each one (that is, for each regular expression). The **XML** sample files ...

Cited by 3 - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

Attribute Grammars for Scalable Query Processing on **XML** Streams - group of 15 »

C Koch, S Scherzinger - Proc. DBPL 2003, 2003 - Springer ... Attribute Grammars for Scalable Query Processing on **XML** Streams ... languages are precisely those recognizable by the deterministic pushdown **automata** (DPDA, see ...

Cited by 9 - [Related Articles](#) - [Web Search](#) - [BL Direct](#)

Regular expressions with numerical occurrence indicators—
preliminary results - group of 2 »

P Kilpelainen, R Tuhkanen - Proc. of the Eighth Symposium
on Programming Languages and ..., 2003 - cs.uku.fi
... way that the SGML standard [13] and the XML and XML
Schema recommendations ... that it
is probably difficult to realize #REs as deterministic
automata, which is ...

Cited by 3 - Related Articles - View as HTML - Web Search

The Design of RELAX NG - group of 2 »

J Clark - 2001 - thaiopensource.com
... is removed. The classic implementation technique for
SGML and XML content models
is to construct a **Glushkov automaton**. The 1-unambiguity ...

Cited by 7 - Related Articles - Cached - Web Search

Validation of XML Document Updates Based on XML
Schema in XML Databases - group of 3 »

SK Kim, M Lee, KC Lee - LECTURE NOTES IN
COMPUTER SCIENCE, 2003 - Springer
... DTD inference for views of XML data. ... From regular
expressions to deterministic **automata**. ...
Communications of the ACM, 11:419–422 (1968) [15]
VM.**Glushkov**. ...

Related Articles - Web Search - BL Direct

Google ►

Result Page: 1 2 3 Next

glushkov automata + xml

Search

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2006 Google



[Home](#) | [Login](#) | [Logout](#) | [?](#)

AbstractPlus

[View TOC](#)

Access this document

 Full Text: [PDF](#) (968 KB)

Download this citation

Choose Citation & Abstract

Download ASCII Text

[Download](#)

» [Learn More](#)

Rights and Permissions

» [Learn More](#)

**BROWSE SEARCH IEEE XP
GUIDE**

A highly-extensible, XML-based language

Dashofy, E.M. van der Hoek, A.
Dept. of Inf. & Comput. Sci., Calif.

This paper appears in: **Software A
Proceedings, Working IEEE/IFL**

Publication Date: 2001

On page(s): 103-112

Meeting Date: 08/28/2001 - 08/31/

Location: Amsterdam, Netherlands

ISBN: 0-7695-1360-3

References Cited: 24

INSPEC Accession Number: 7092

Digital Object Identifier: 10.1109/

Posted online: 2002-08-07 00:11:0

Abstract

Software architecture research focus architectures as specified in archite (ADLs). As research progresses in architectures, more and more archi created. Ideally, this information ca An extensible modeling language i with and building tools for novel m from evolving research. Traditiona small set of modeling constructs ve poorly. XML provides an ideal pla an extensible modeling language fc Previous XML-based ADLs succes large base of off-the-shelf tool sup advantage of its extensibility. To g researchers more freedom to exploi modeling techniques, while maxim modeling constructs, we have deve extensible XML-based ADL. xAD design time modeling, architecture, and model-based system instantiati has a set of extensible infrastruc creation, manipulation, and sharing

Index Terms
Insne

Controlled Indexing

[hypermedia markup languages](#)
[architecture](#) [software reusabi](#)
[languages](#)

Non-controlled Indexing

ADLs XML-based ADLs ar
information architecture confi
architecture description langua
modeling extensible infrastru
modeling language highly ex
architecture description langua
system instantiation modelin
techniques off-the-shelf tool
architecture specification soft

Author Keywords

Not Available

References

No references available on IEEE X

Citing Documents

No citing documents available on I

◀ [View TOC](#) | [Back to Top](#) ▶

Indexed by
 Inspec

E

©



[Subscribe \(Full Service\)](#) [Register \(Limited Service\)](#)

Search: The ACM Digital Library The ACM

THE GUIDE TO COMPUTING LITERATURE

[Feedback](#) [Report a problem](#)

Handling syntactic constraints in a DTD-compliant XML editor

Full text [Pdf \(364 KB\)](#)

Source [Document Engineering archive](#)
Proceedings of the 2003 ACM symposium on Document engine
Grenoble, France
SESSION: Editing and authoring [table of contents](#)
Pages: 222 - 224
Year of Publication: 2003
ISBN:1-58113-724-9

Authors [Y. S. Kuo](#) Academia Sinica, Taiwan
[Jasper Wang](#) Academia Sinica, Taiwan
[N. C. Shih](#) Academia Sinica, Taiwan

Sponsors [SIGWEB](#): ACM Special Interest Group on Hypertext, Hypermedia
[SIGIR](#): ACM Special Interest Group on Information Retrieval
[ACM](#): Association for Computing Machinery

Publisher ACM Press New York, NY, USA

Additional Information: [abstract](#) [references](#) [citations](#) [index terms](#) [collaborative peer](#)

Tools and Actions: [Find similar Articles](#) [Review this Article](#)
[Save this Article to a Binder](#) Display Formats: [BibTeX](#)

DOI Bookmark: Use this link to bookmark this Article: <http://doi.acm.org/10.1145/1055558.1055559>
[What is a DOI?](#)

↑ ABSTRACT

By exploiting the theories of automata and graphs, we propose algorithms and valid XML documents [4][5]. The editing process avoids syntactic violations the user from any syntactic concerns. Based on the proposed algorithms and p

XML editor with forms as its user interface.

↑ REFERENCES

Note: OCR errors may be found in this Reference List extracted from the full opted to expose the complete List rather than only correct and linked reference

- 1 Anne Brüggemann-Klein, Regular expressions into finite automata, Theor. Science, v.120 n.2, p.197-213, Nov. 22, 1993
- 2 D. D. Cowan , E. W. Mackie , G. M. Pianosi , G. de V. Smit, Rita—an editor for manipulating structured documents, Electronic Publishing—Origination, I Design, v.4 n.3, p.125-150, Sept. 1991
- 3 Shimon Even, Graph Algorithms, W. H. Freeman & Co., New York, NY,
- 4 W3C, Extensible Markup Language (XML) 1.0, W3C Rec., Feb. 10, 1998
- 5 W3C, XML Schema Part 1: Structures, W3C Rec., May 2, 2001.
- 6 XMLSoftware, <http://www.xmlsoftware.com/editors.html>

↑ CITINGS 2

Y. S. Kuo , N. C. Shih , Lendle Tseng , Hsun-Cheng Hu, Generating form-based XML vocabularies, Proceedings of the 2005 ACM symposium on Document 02-04, 2005, Bristol, United Kingdom

Marc Dymetman, Chart-parsing techniques and the prediction of valid editing document authoring, Proceedings of the 2004 ACM symposium on Document 28-30, 2004, Milwaukee, Wisconsin, USA

↑ INDEX TERMS

Primary Classification:

H. Information Systems

↳ H.5 INFORMATION INTERFACES AND PRESENTATION (I.7)

↳ H.5.2 User Interfaces (D.2.2, H.1.2, I.3.6)

↪ **Subjects:** Theory and methods

Additional Classification:

H. Information Systems

↪ H.5 INFORMATION INTERFACES AND PRESENTATION (I.7)

↪ H.5.2 User Interfaces (D.2.2, H.1.2, I.3.6)

↪ **Subjects:** Graphical user interfaces (GUI); Interaction styles (e.g., c/forms, direct manipulation)

General Terms:

Algorithms, Design

Keywords:

XML editor, automata theory, regular expression

↑ **Collaborative Colleagues:**

Y. S. Kuo:	C. Chen	F. Ruskey	S. Wu
	T. C. Chern	N. C. Shih	
	W. K. Chou	W.-K. Shih	
	Tyng-Ruey Chuang	Wei-kuan Shih	
	H. F. Hu	M. T. Shing	
	Hsun-Cheng Hu	S. J. Su	
	T. C. Hu	J. C. Tsay	
	S. -Y. Hwang	Lendle Tseng	
	Chuan-Chieh Jung	Chien-Min Wang	
	Wen-Min Kuan	Jaspher Wang	
N. C. Shih:	Hsun-Cheng Hu		
	Y. S. Kuo		
	Lendle Tseng		
	Jaspher Wang		
Jaspher Wang:	Y. S. Kuo		
	N. C. Shih		

↑ **Peer to Peer - Readers of this Article have also read:**




- Data structures for quadtree approximation and compression **Communica** 28, 9

Hanan Samet

- A hierarchical single-key-lock access control using the Chinese remainder of the 1992 ACM/SIGAPP Symposium on Applied computing
Kim S. Lee , Huizhu Lu , D. D. Fisher
- The GemStone object database management system **Communications of**
Paul Butterworth , Allen Otis , Jacob Stein
- Putting innovation to work: adoption strategies for multimedia communication **Communications of the ACM** 34, 12
Ellen Francik , Susan Ehrlich Rudman , Donna Cooper , Stephen Levine
- An intelligent component database for behavioral synthesis **Proceedings of ACM/IEEE conference on Design automation**
Gwo-Dong Chen , Daniel D. Gajski

The ACM Portal is published by the Association for Computing Machinery
ACM, Inc.

[Terms of Usage](#) [Privacy Policy](#) [Code of Ethics](#) [Contact](#)

Useful downloads:  [Adobe Acrobat](#)  [QuickTime](#)  [Windows Media Player](#)

**Rita an Editor and User Interface
for Manipulating Structured**

Documents (1991) ([Make
Corrections](#)) ([18 citations](#))

D.D. Cowan, E.W. MacKie, G.M.
PIANOSI, G. de V. Smit
Electronic Publishing

CiteSeer
Electronic Literature Digital Library

[Home/Search](#)

[Bookmark](#) [Context](#) [Related](#)

View or download:

cajun.cs.nott.ac.uk/wiley...ep048dc.pdf

Cached: [PS.gz](#) [PS](#) [PDF](#)

[Image](#) [Update](#) [Help](#)

From: cajun.cs.nott.ac.uk/wi...epoddkwi
([more](#))

([Enter author homepages](#))

([Enter summary](#))

Rate this article: 1 2 3 4 5 (best)

[Comment on this article](#)

Abstract: This paper describes Rita, its user interface and some of its internal structure and algorithms, and relates anecdotal user experiences. Comparisons are also made with other commercial and experimental systems. ([Update](#))

Context of citations to this paper: [More](#)

...then all the visualization information that can be generated from the tags will not be present. **Syntax directed editors such as Rita [CMPS91] the SoftQuad Author Editor [Sof89] and editors produced by the Cornell Synthesizer [RT87] could be used to address some of these...**

...a WYSIWYG manner. The syntax of SGML looks as cumbersome and difficult to read as typical WEB. **With some SGML viewers and structure editors [2] the user is not even aware that the underlying document is tagged with SGML.** An SGML markup language and accompanying style sheet can...

Cited by: [More](#)

A Structural Adviser for the XML Document Authoring - Chidlovskii (2003)
([Correct](#))

Electronic Publishing, Vol. 8(2 3), 125--138 (june.. - New Presentation
Language ([Correct](#)))

CIRL/PIWI: A GUI Toolkit Supporting Retargetability - Cowan, Durance..

(1993) ([Correct](#))

Active bibliography (related documents): [More](#) [All](#)

1.2: Incremental Updates in Structured Documents - Lindén (1994)
([Correct](#))

0.9: Interactively Editing Structured Documents - Furuta, QUINT, André (1988) ([Correct](#))

0.6: Transition Diagram Systems and Normal Form Algorithms - Giammarresi, Wood (1998) ([Correct](#))

Similar documents based on text: [More](#) [All](#)

0.2: RITA: Receiver Initiated Just-in-Time Tree Adaptation for .. - Xu, Tang, Banerjee, Lee (2003) ([Correct](#))

0.2: Style Control in the Quill Document Editing System - Wolfsthal (1991) ([Correct](#))

0.2: Practical Language-Based Editing For Software Engineers - Vanter (1995) ([Correct](#))

Related documents from co-citation: [More](#) [All](#)

7: Oxford University Press (context) - Goldfarb, Handbook - 1992

4: Reading source code (context) - Raymond - 1991

4: Free Software Foundation (context) - Stallman, Manual - 1994

BibTeX entry: ([Update](#))

D. Cowan, E. Mackie, G. Pianosi, and G. d. V. Smit, "Rita -- An Editor and User Interface for Manipulating Structured Documents," Electronic Publishing, Origination, Dissemination and Design, vol. 4, pp. 125--150, September 1991. <http://citeseer.ist.psu.edu/cowan91rita.html> [More](#)

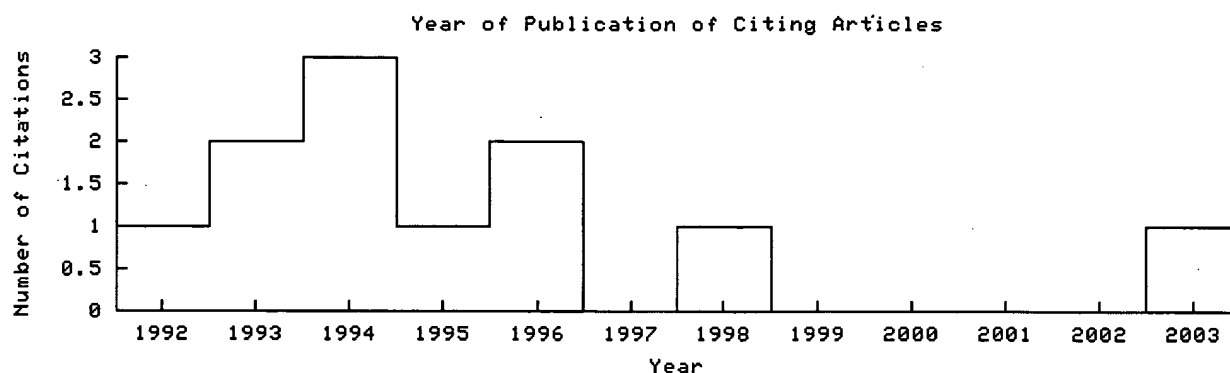
```
@article{ cowan91rita,  
  author = "Donald D. Cowan and E. W. Mackie and G. M. Pia  
  title = "Rita - an Editor and User Interface for Manipul  
  journal = "Electronic Publishing",  
  volume = "4",  
  number = "3",  
  pages = "125-150",  
  year = "1991",  
  url = "citeseer.ist.psu.edu/cowan91rita.html" }
```

Citations (may not include all citations):

- 1894 Introduction to Automata Theory (context) - Hopcroft, Ullman - 1979
- 226 A Document Preparation System (context) - Lamport - 1986
- 169 The Synthesizer Generator: A System for Constructing Language.. (context) - Reps, Teitelbaum - 1987
- 70 The Cornell Program Synthesizer: a syntax-directed programming.. (context) - Teitelbaum, Reps - 1981
- 70 The synthesizer generator (context) - Reps, Teitelbaum - 1984
- 59 Theory of Computation (context) - Wood - 1987
- 47 Interactively editing structured documents - Furuta, Quint et al. - 1988
- 47 Grif: an interactive system for structured document manipulation.. (context) - Quint, Vatton
- 27 The Systems Programming Series (context) - Foley, van Dam et al. - 1990
- 20 Document Style Semantics and Specification Language (context) - for, New - 1991
- 16 Specifying structured document transformations (context) - Furuta, Stotts
- 13 Multiple representation document development (context) - Chen, Harrison - 1988
- 7 Microsoft Word (context) - Corporation, Washington - 1989
- 6 Standard for Electronic Manuscript Preparation and Markup (context) - American, Electronic - 1987
- 6 Defining document styles for WYSIWYG processing (context) - Chamberlin, Hasselmeier et al.
- 5 Syntax-directed editing: towards integrated programming environment.. (context) - Medina-Mora - 1982
- 4 Formatting structured documents: batch versus interactive (context) - Coray, Ingold et al.
- 4 Editor User's Manual (context) - Inc, Canada et al. - 1989
- 4 Quill: An extensible system for editing documents of mixed type.. (context) - Chamberlin, Hasselmeier et al. - 1988
- 4 Supporting document development with Concordia (context) - Walker - 1988
- 4 An interactive prototyping environment for language design (context) - Feiler, Fahimeh et al. - 1986
- 4 A formalization of transition diagram systems (context) - Lomet - 1973
- 3 An integrated, but not exact-representation, editor/formatter.. (context) - Furuta
- 3 NROFF/TROFF user (context) - Ossana - 1976

- 3 A Formatter-Independent Structured Document Preparation Syst.. (context)
- de and, Smit - 1987
- 3 The SAGA project: a system for software development (context) -
Campbell, Kirsulis - 1984
- 2 A structure editor for abstract document objects (context) - Kimura - 1986
- 2 Manipulating partial documents in a syntax directed environm.. (context) -
de and, Smit et al. - 1990
- 2 Incremental execution environment (context) - Bhatti - 1988
- 2 Is what you see enough to get? A description of the Interlea.. (context) -
Morris - 1985
- 1 Display-oriented structure manipulation in a multi-purpose s.. (context) -
Feiler, Kaiser - 1983
- 1 Combining interactive document editing with batch document f.. (context)
- de and, Smit et al.
- 1 SYNED -- a language-based editor for an interactive programm.. (context)
- Gansner - 1983
- 1 Experiences with RITA (context) - Cowan, Mackie et al. - 1990
- 1 Text Processing and Document Manipulation (context) - Int - 1986
- 1 WATCOM Publications Limited (context) - Mackie, Pianosi et al. - 1991
- 1 Cambridge University Press (context) - Manipulation, ed et al. - 1988
- 1 The implementation of Etude, an integrated and interactive d.. (context) -
Hammer - 1981
- 1 Document convergence in an interactive formatting system (context) -
Chamberlin - 1987
- 1 Motif User's Guide (context) - Foundation - 1990
- 1 Document Composition Facility: General Markup Language Start..
(context) - White, New - 1989
- 1 of Waterloo, Waterloo, Ontario, Canada, Waterloo SCRIPT GML ..
(context) - Computing, University - 1988
- 1 PEN: A hierarchical document editor (context) - Allen, Nix et al. - 1981
- 1 WATCOM Publications Limited (context) - McKee, Welch et al. - 1990
- 1 An overview of the W document preparation system (context) - King
- 1 WRITE-IT SGML Editor (context) - Systems, UK
- 1 RITA---MANIPULATING STRUCTURED DOCUMENTS (context) -
Goldfarb, Handbook et al. - 1990
- 1 SGML: A standard language for text description (context) - van Huu -
1985
- 1 Waterloo Rita Document Class Generator --- Reference (context) - Pianosi

- 1990



The graph only includes citing articles where the year of publication is known.

Documents on the same site

(<http://cajun.cs.nott.ac.uk/wiley/journals/epobetan/pdf/epoddkwi.html>):
More

Journal publishing with Acrobat: the CAJUN project - Philip Smith (1993)
(Correct)

Distributed Documents: an architecture for open.. - Hatzimanikatis..
(Correct)

Tools for Printing Indexes - Bentley, Kernighan (1988) (Correct)

Online articles have much greater impact More about CiteSeer.IST Add
search form to your site Submit documents Feedback

CiteSeer.IST - Copyright Penn State and NEC

TYPE OF PROPOSAL: paper.

TITLE: Automatic DTD simplification by examples

KEYWORDS: digital libraries, markup, automatic learning.

p. 5

Add

eliminate

AUTHOR: Alejandro Bia

AFFILIATION: Miguel de Cervantes Digital Library, University of Alicante

E-MAIL: abia@dlsi.ua.es

AUTHOR: Rafael C. Carrasco

AFFILIATION: Dept. Lenguajes y Sistamas Inform'ticos,
University of Alicante

E-MAIL: carrasco@dlsi.ua.es

CONTACT ADDRESS: Biblioteca Virtual Miguel de Cervantes.Universidad de Alicante. Apdo. Correos 99, 03080, Alicante, Spain.

FAX NUMBER: 34-96-590-9477

PHONE NUMBER: 34-96-590-9567

Automatic DTD simplification by examples

This paper describes a method for the automatic generation of simplified DTDs from a source DTD and a set of sample marked up files. The purpose is to create the minimum DTD that the sample set of files comply. In this way, new files can be created and parsed using this simplified DTD but still being compliant to the original, more general DTD. The simplified DTD can be used to make the task of markup easier, specially for non-experienced XML writers.

The resulting tool was used at the Miguel de Cervantes digital library (<http://cervantesvirtual.com/>) to obtain simplified versions of the TEI.DTD (Sperberg-McQueen and Burnard, 1994). This work is part of a larger project in the field of text markup and derived applications (Bia and Pedreño, 2000).

Motivation

"Having standardized-XML-vocabularies for common things allows

developers to reuse existing DTDs, saving the cost of developing custom DTDs. Custom DTDs isolate their users and applications from others that might otherwise be able to share commonly formatted documents and data. Shared DTDs are the foundation of XML data interchange and reuse" (Hunter, 2000).

Saving the cost of developing our own DTD, and text interchangeability are some of the reasons why the **teixlite.dtd** (XML version of the SGML **teilight.dtd** of the TEI encoding scheme) has been chosen at the Cervantes digital library, but the **TEIxlite** is still too complex for markup beginners. Our markup team is composed mostly of humanists with some computer skills but who appreciate their computer work be simplified as much as possible.

On the other hand our XML documents do not use, and do not need all the markup options provided by the **teixlite.dtd**. So a simpler DTD was needed to simplify markup tasks and to avoid possible use of unwanted markup options. But we still wanted our files to be TEI compliant and benefit from the advantages of sharing a common DTD with other international digitization projects. In brief, we needed a simpler DTD, a TEI compliant DTD, that is a valid subset of the **teilight.dtd**.

We started by defining the kinds of modifications we will allow ourselves to make to the TEIlight DTD, in order to make it simpler to use but at the same time keeping our documents TEI-compatible (except for minor exceptions). In this sense we allowed the following changes:

- To add normalized values to some attributes in order to force the use of fixed values instead of free data entry.
- To add new attributes only in a few necessary cases (this is the only exception that may keep our files from being TEI compliant, but we thought that these added attributes can be easily eliminated anytime we wanted to comply the TEI standard).
- To make restrictions in element inclusion rules (we wanted to eliminate the possibility of including certain elements at certain levels of the markup).
- To make some optional elements/attributes mandatory to force following our specific markup norms.
- To eliminate optional elements we will not use to simplify the markup task and to avoid possible errors (basically we wanted to eliminate the

features we decided not to use)

It is clear that doing the simplifications by hand is tedious and error prone. Constructing a set of sample documents representative of all the types of documents we need to markup together with a program that simplifies the DTD automatically will alleviate this task.

Previous works

Document types are defined by extended context-free grammars where the right hand side of productions are unambiguous regular expressions (Bruggemann, 1998). Previous work has addressed the task of identifying a DTD from examples. A common difficulty in this approach is the need to find a correct degree of generalization. Some practical tools as FRED (Shafer, 1995) let the users customize their preferred degree of generalization. Ahonen builds a (k,h)-testable model (Ahonen, 1995; Ahonen, 1997; Ahonen, Mannila, and Nikunen, 1997).

Young-Lai and Tompa (Young-Lai and Tompa, 2000) rely on a stochastic approach to control overgeneralization, based in turn on the algorithm by Carrasco and Oncina (Carrasco, 1998). Presumably, the stochastic approach needs large collections of hand-tagged documents.

Pizza-Chef (Burnard, 1997) is a tool to generate TEI-compliant DTDs suited to a particular task. In this case, predefined tasks and TEI DTDs are only allowed.

Objectives

However, a general DTD defining a global frame that a whole set of files must fulfill allows for a natural way to avoid overgeneralization. In this sense, any particularized, narrow scope DTD should not accept any document that is not accepted by the general, wide scope DTD.

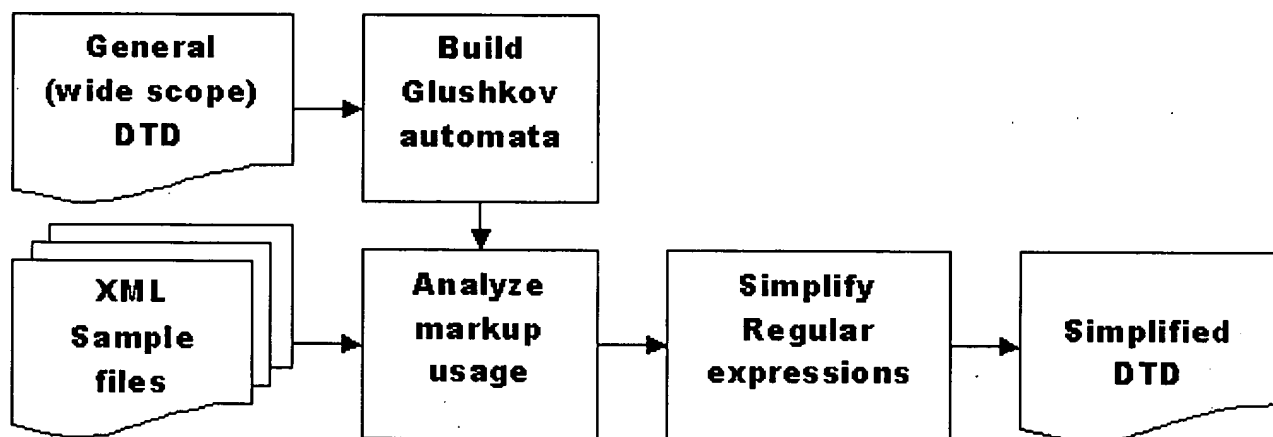
Therefore, the objective of our approach is to automatically select only those DTD features that are used by a set of valid documents (validated against the more general DTD) and eliminate the rest of them, obtaining a narrow scope DTD which defines a subset of the original markup scheme. This "pruned" DTD can be used to build new documents of the same markup subclass, which

in turn would still comply the original general DTD. Needless to say that working with a simpler DTD is easier.

General description

For the implementation of the DTDprune toolkit we needed both an XML and a DTD parser. We assumed that both the XML sample files and the source DTD would be well-formed and valid, so there would be no need to build validating parsers. Instead, we developed two simple parsers, based on the XML BNF Grammar described in (Harold, 1999). A diagram of the process is shown below in the figure.

Architecture of the DTD simplifier



As the diagram shows, the general DTD is processed to extract the structure of the markup model with which we build a Glushkov automata (Caron and Ziadi, 2000). The XML sample files are preprocessed to extract the elements used and their nesting patterns. Based on the Glushkov automata that represent the regular expressions that define the possible element contents according to the general DTD, we keep track of the elements used in the sample files and mark the visited states of the automata. Finally, a simplification process takes place. This process eliminates unused elements and simplifies the right parts of element definitions, i.e. the regular expressions that define further nestings. The simplified DTD structure is used to generate the new simplified DTD.

Conclusions

Using this automated method, the simplified DTD can be updated

immediately in the event that new features are added to (or eliminated from) the sample set of XML files (modifications to files of the sample-set must be done using the general DTD for validation). This process can be repeated to incrementally produce a final narrow-scope DTD. In this way, we use a complex DTD as a general markup-design frame to build a simpler working-DTD that suits a specific project's markup needs.

Another use of this technique is to build a one-document DTD, i.e. the minimum DTD derived from the general DTD that a given XML document would comply.

Another benefit of this technique is that we can produce statistics that may help markup designers improve their markup schemes. Information about the frequency of use of certain elements within others, helps us to detect unusual structures that could reflect mark-up mistakes, misuse of the DTD, or DTD features that may allow unwanted generalization. This statistical data on the use of markup may help us take decisions about adding new markup constraints, or on the contrary expand the simplified DTD.

References

C.M. Sperberg-McQueen and Lou Burnard, editors.

Guidelines for Electronic Text Encoding and Interchange (Text Encoding Initiative P3), Revised Reprint, Oxford, May 1999.

TEI P3 Text Encoding Initiative, Chicago - Oxford, May 1994.

Alejandro Bia and Andrés Pedreño.

The Miguel de Cervantes Digital Library: The Hispanic Voice on the WEB. *LLC (Literary and Linguistic Computing) journal, Oxford University Press*, (to be published soon) 2000.

Presented at ALLC/ACH 2000, The Joint International Conference of the Association for Literary and Linguistic Computing and the Association for Computers and the humanities, 21/25 July 2000, University of Glasgow.

Anne Brüggemann-Klein and Derick Wood.

One-unambiguous regular languages.

Information and Computation, 142(2):182--206, 1 May 1998.

Keith E. Shafer.

Creating dtDs via the gb-engine and fred.

Technical report, OCLC Online Computer Library Center, Inc., 6565 Frantz Road, Dublin, Ohio 43017-3395, 1995.

Helena Ahonen.

Automatic generation of SGML content models.

Electronic Publishing Origination, Dissemination, and Design, 8(2/3):195--206, June\September 1995.

H.Ahonen, H.Mannila, and E.Nikunen.

Generating grammars for SGML tagged texts lacking DTD.

Mathematical and Computer Modelling, 26(1):1--13, 1997.

H.Ahonen.

Disambiguation of SGML content models.

Lecture Notes in Computer Science, 1293:27, 1997.

Matthew Young-Lai and Frank W. M. Tompa.

Stochastic grammatical inference of text database structure.

Machine Learning, 40(2):1, 2000.

Rafael C. Carrasco and Jose Oncina.

Learning deterministic regular grammars from stochastic samples in polynomial time.

RAIRO (Theoretical Informatics and Applications), 33(1):1--20, 1999.

Lou Burnard.

The Pizza Chef: a TEI Tag Set Selector.

<http://www.hcu.ox.ac.uk/TEI/pizza.html>, September 1997.

(Original version 13 September 1997, updated July 1998; Version 2 released 8 Oct 1999).

David Hunter, Curt Cagle, Dave Gibbons, Nikola Ozu, Jon Pinnock, and Paul Spencer.

Beginning XML.

Programmer to Programmer. Wrox Press, 1102 Warwick Road, Acocks Green, Birmingham, B27 6BH, UK, 1st edition, 2000.

Pascal Caron and Djelloul Ziadi.

Characterization of Glushkov automata.

TCS: Theoretical Computer Science, 233:75--90, 2000.

Robert D. Cameron.

REX: XML Shallow Parsing with Regular Expressions.

Markup Languages: Theory & Practice, 1(3):61--68, Summer 1999.



glushkov automata

1980

- 2

Scholar [All articles](#) [Recent articles](#) Results 1 - 10 of about 339 for **glushkov**

All Results[A Brüggemann-K...](#)[D Wood](#)[P Caron](#)[A Brüggemann-K...](#)[D Ziadi](#)

[PS] [Characterization of Glushkov automata - group of 6 »](#)

P Caron, D Ziadi - TCS, 2000 - dlsi.ua.es

Page 1. Characterization of **Glushkov automata** ...

1 Page 2. structural properties of

Glushkov automata and to establish a characterization of these **automata**. ...

[Cited by 31](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

→ **[PS]** [Regular expressions into finite automata - group of 4 »](#)

→ A Brüggemann-Klein - Theoretical Computer Science, 1993 - informatik.tu-muenchen.de ... deterministic regular expressions, ie expressions whose **Glushkov automaton** is deterministic, as a description language for document types. ...

[Cited by 78](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

[One-Unambiguous Regular Languages - group of 15 »](#)

A Brüggemann-Klein, D Wood - Information and Computation, 1998 - Elsevier

... Berry and Sethi [BS86] showed that **Glushkov automata** are natural representations of regular expressions. Definition 2.3. We define the **Glushkov automaton** G ...

[Cited by 83](#) - [Related Articles](#) - [Web Search](#) - [BL Direct](#)

[PS] [The Validation of SGML Content Models - group of 11 »](#)

A Brüggemann-Klein, D Wood - MATH COMPUT MODELL(OXFORD), 1997 - informatik.tu-muenchen.de

... We also define deterministic regular expressions based on **Glushkov automata** and discuss the determinism of content models. ... a 1 to a 4 in the **Glushkov automaton**. ...

[Cited by 25](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#) - [Library Search](#) - [BL Direct](#)

[Glushkov construction for multiplicities - group of 4 »](#)

P Caron, M Flouret - Pre-Proceedings of CIAA, 2000 - Springer

... 3 The Extended **Glushkov Automaton** ... All these functions allow us to define an extended **Glushkov automaton**. Page 6. 72 P. Caron and M. Flouret ...

[Cited by 7](#) - [Related Articles](#) - [Web Search](#) - [BL Direct](#)

[PS] [Deterministic regular languages - group of 5 »](#)

A Brüggemann-Klein, D Wood - STACS, 1992 - informatik.uni-freiburg.de

... E the **Glushkov automaton** of E. Figure 1 shows the two **Glushkov automata** corresponding to the expressions ... Figure 1 The **Glushkov automata** for $(a + b)$...

[Cited by 16](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

[Determinization of Glushkov automata - group of 3 »](#)

JM Champarnaud, D Ziadi, JL Ponty - Third International Workshop on Implementing Automata-WIA' ..., 1999 - Springer

Page 1. Determinization of **Glushkov Automata** ...
Section 3 describes the subset construction.
Section 4 gathers our results about **Glushkov automata** determinization. ...

[Cited by 4](#) - [Related Articles](#) - [Web Search](#) - [BL Direct](#)

[An Optimal Parallel Algorithm to Convert a Regular Expression into its **Glushkov Automaton** - group of 7 »](#)

D Ziadi, JM Champarnaud - TCS, 1999 - jmc.feydakin.org

... The aim of this paper is to describe a CREW-PRAM optimal algorithm which converts a regular expression of size s into its **Glushkov automaton** in $O(\log s)$ time ...

[Cited by 3](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#)

[BOOK] [Unambiguous Regular Expressions and SGML Document Grammars](#)

A Brüggemann-Klein, D Wood - 1992 - csd.uwo.ca

... Berry and Sethi [BS86] show that **Glushkov automata** are natural representations of regular expressions. Denition 2.3 We dene the **Glushkov automaton** M ...

[Cited by 9](#) - [Related Articles](#) - [View as HTML](#) - [Web Search](#) - [Library Search](#)

[An efficient null-free procedure for deciding regular language membership - group of 5 »](#)

JL Ponty - Theoretical Computer Science, 2000 - ingentaconnect.com

... It is based on our ZPC representation of the **Glushkov automaton** of a regular expression. This procedure requires a specific representation ...

[Cited by 5](#) - [Related Articles](#) - [Web Search](#) - [BL](#)

[Direct](#)

Goooooooooooooogle ►

Result Page: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [Next](#)

glushkov automata

Search

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2006 Google